

Faculty of Computer Science, Dalhousie University

7-Oct-2023

CSCI 4152/6509 — Natural Language Processing

Lecture 9: P0 Topics Discussion; Introduction to Probabilistic Modeling

Location: Carleton Tupper Building Theatre C Instructor: Vlado Keselj
Time: 16:05 – 17:25

Previous Lecture

- CNG classification method
 - Edit distance:
 - introduction, properties, dynamic programming approach, example, algorithm
-

P0 Topics Discussion

- Discussion of individual projects as proposed in P0 submissions
- Projects discussed: P-01, P-03, P-04, P-05, P-06

Part III

Probabilistic Approach to NLP

11 Introduction to Probabilistic Approach to NLP

11.1 Logical versus Plausible Reasoning

Artificial Intelligence (AI) is an area of Computer Science related to the task of developing software and computers that exhibit some form of intelligent behaviour. NLP can be considered to be a sub-area of AI. Solving problems in AI generally involves implementing in a computer some form of reasoning or inference. The automated inference methods can be divided into two large groups: *logical reasoning*, and *plausible reasoning*. As we will see, this division can be adopted for the NLP area as well.

1. Logical reasoning is known also as classical, symbolic, or knowledge-based AI. The rule-based AI is also a form of logical reasoning. This form of automated inference is based on the principles that we can find in mathematical logic. We start with basic set of premises, known as axioms, and using a fixed set of rules, we derive new conclusions. This type of reasoning is called *monotonic* since with more evidence, i.e., more known facts, we can produce more conclusions. A consequence of this is that once we make conclusions, they cannot be retracted or canceled. In other words, if we conclude that something is true, we cannot conclude that it is wrong based on the new evidence because that would mean that we either had some wrong evidence in the first place, or we made a mistake in reasoning. We also describe this reasoning as *certain*, since the model is designed in such way that we can derive only certain (or definite, guaranteed) conclusions.

2. Plausible reasoning is an automated reasoning in which we typically derive plausible conclusions; i.e., conclusions that may be true, and we usually have some way of rating the truthfulness of such conclusions. There are different approaches to this kind of automated reasoning. The most widely used are the approaches based on

the mathematical probability theory, where we try to make a mathematically sound estimate of the probability that our conclusions are true. Other approaches include fuzzy logic and neural networks. Unlike logical reasoning, plausible reasoning is *uncertain*, since the conclusions are not guaranteed to be correct. Logical reasoning is also *non-monotonic* since based on new evidence, we can withdraw some conclusions, and with more evidence we may end up with less conclusions. Plausible reasoning allows for contradicting evidence.

Plausible Reasoning

- How to combine ambiguous, incomplete, and contradicting evidence to draw reasonable conclusions?
- Frequently approached as the task of making plausible inference of some hidden structure from observations.
- examples:

Input (observations)		Hidden Structure
symptoms	→	illness
pixel matrix	→	object, relations
speech signal	→	phonemes, words
word sequence	→	meaning
sentence	→	parse tree
word sequence	→	POS tags, names, entities
words in e-mail Subject:	→	Is message spam? Yes/No
text	→	text category (class)

Probabilistic NLP as a Plausible Reasoning Approach

- Regular expressions and finite automata are example of logical or knowledge-based approach to NLP
- Plausible approaches to NLP:
 1. Probabilistic: use of Theory of Probability, also known as stochastic or statistical NLP
 - Alternative plausible approaches, examples:
 2. neural networks,
 3. kernel methods,
 4. fuzzy logic, fuzzy sets,
 5. Dempster-Shafer theory
 6. rough sets,
 7. default logic, ...

11.2 Review of Basics of Probability Theory

This is a brief and intuitive review of some basic notions from the theory of probability.

- You should have this background from previous courses; this is just a review,
 - discussed a bit in the textbook: [JM] 5.5, and [MS] 2.1
- Simple event or basic outcome
 - e.g., rolling a die, choosing a letter
- *Event space*: the set of all outcomes, usually denoted Ω

The event space is also known as the sample space. For example, the event space for rolling a die has six elements, corresponding to the six different outcomes; or an event space for choosing a letter of the English alphabet has 26 elements if we ignore letter case.

- *Event or outcome* is a set of simple events or basic outcomes
- In other words event is any subset of Ω ; i.e., $A \subseteq \Omega$
- Each event is associated with a probability, which is a number between 0 and 1, inclusive: $0 \leq P(A) \leq 1$

Probability Examples

- P(“rolling a 6 with a die”) = $1/6$
- Choosing a letter of English alphabet:
 - If we choose uniformly: $P('a') = 1/26 \approx 0.04$
 - Choosing from a text: $P('a') \approx 0.08$
 - Remember our output from “Tom Sawyer”:

```
35697 0.1204 e
28897 0.0974 t
23528 0.0793 a
23264 0.0784 o
20200 0.0681 n
...
```

As we saw from the “Tom Sawyer” example, the probability of the letter ‘a’ in a typical English text is about 0.08.

The probability can be seen as a function that maps a set of possible experiment outcomes to a real number between 0 and 1, including 0 and 1, i.e. to a number from the interval $[0, 1]$. We always have in mind certain space of outcomes Ω and we assume that in each experiment (trial, instance, or model configuration), one of those outcomes will happen. The probability that one of the outcomes from a set A ($A \subset \Omega$) will happen, is denoted $P(A)$. A set of outcomes A is called an event.

Probability Axioms

- Probability axioms: nonnegativity, additivity, normalization:

The basic properties of the probability, known as the probability axioms, are:

- **(Nonnegativity)** $P(A) \geq 0$, for any event A
- **(Additivity)** for disjoint events A and B , i.e., if $A, B \subset \Omega$ and $A \cap B = \emptyset$, then

$$P(A \cup B) = P(A) + P(B).$$

More generally, for a possibly infinite sequence of disjoint events A_1, A_2, \dots ,

$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$$

- **(Normalization)** $P(\Omega) = 1$, where Ω is the entire sample space.
- Some consequences of the above axioms are: $P(\emptyset) = 0$ and $P(\Omega - A) = 1 - P(A)$

Independent and Dependent Events

- Independent events A and B (definition): $P(A, B) = P(A) \cdot P(B)$
- Use of comma in: $P(A, B) = P(A \cap B)$
- Example: choosing two letters in text
 1. Choosing independently: $P('t') = 0.1, P('h') = 0.07, P('t', 'h') = 0.007$
 2. Choosing two consecutive letters (dependent events): $P('t', 'h') = 0.04$

Conditional Probability

- Conditional probability

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

The probability $P(A|B)$ (it is read “probability of A given B”) is the probability of event A happening given that the event B happens.

- Expressing independency using conditional probability
Two events A and B are independent if and only if:

$$P(A|B) = P(A)$$

This is an alternative definition of independent events.

Annotation with More Events

- There is a bit of flexibility in using notation; e.g.,
- $P(A, B, C) = P(A \cap B \cap C)$
- $P(A|B, C) = P(A|B \cap C)$
- $P(A, B, C|D, E, F) = P(A \cap B \cap C|D \cap E \cap F)$
- and so on.
- Three independent events: $P(A, B, C) = P(A)P(B)P(C)$
- Conditionally independent events

$$P(A, B|C) = P(A|C) \cdot P(B|C)$$

Bayes' Theorem

- Bayes' theorem (one form):

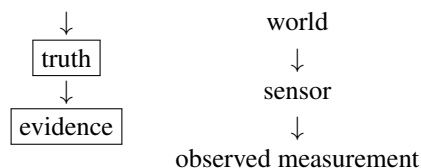
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- The second form is based on breaking the set B into disjoint sets $B = A_1 \cup A_2 \cup \dots$:

$$P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{\sum_i P(B|A_i)P(A_i)}$$

11.3 Bayesian Inference and Generative Models

In our further discussion, we will use *Bayesian Inference* applied to the *Generative Models*. These models are called the generative models because they describe how a particular event that we are studying is generated; i.e., what were the causal relations that were involved in producing certain observations. We can use the following visual representation to describe this model in general:



We assume that there is certain true structure or information that caused some events, which can be observed through some sensors. We will simply refer to this true information as *truth* and information that is available to us as *observable evidence*, or simply as *evidence*.

Notation Remark: 'max' and 'arg max' Operators. Before presenting the main use of the Bayes theorem in Bayesian inference, let us first clarify the use of 'max' and 'arg max' operators in formulae.

The operator 'max' is applied to an expression depending on an index variable x (denoted \max_x) and it returns the maximum value of the expression over all possible values of the variable x . On the other hand, the operator 'arg max' gives one value of the index variable x (denoted $\arg \max_x$) for which the expression achieves the maximum value. Figure 3 illustrates graphically these operators for a function $y = f(x)$.

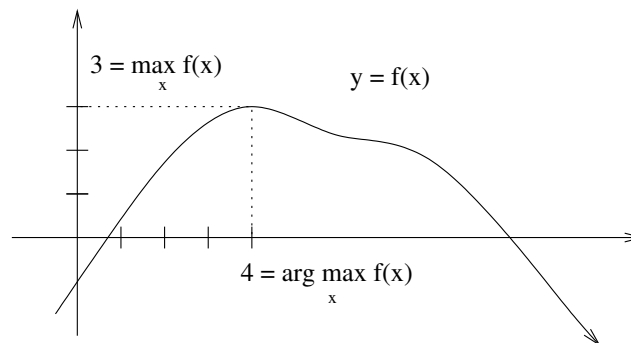


Figure 3: Illustrating operators max and arg max

Bayesian Inference: Using Bayes' Theorem

- Bayesian inference is a principle of combining evidence

$$\begin{aligned} \text{conclusion} &= \arg \max_{\text{possible truth}} P(\text{possible truth}|\text{evidence}) \\ &= \arg \max_{\text{possible truth}} \frac{P(\text{evidence}|\text{possible truth})P(\text{possible truth})}{P(\text{evidence})} \\ &= \arg \max_{\text{possible truth}} P(\text{evidence}|\text{possible truth})P(\text{possible truth}) \end{aligned}$$

- application to speech recognition: acoustic model and language model

The above equations describe Bayesian inference. To understand steps in the equations, let us go over them step by step:

$$\text{conclusion} = \arg \max_{\text{possible truth}} P(\text{possible truth}|\text{evidence})$$

The above equation states that we reach the conclusion about the best possible truth by finding the maximal probability of that ‘truth’ given the evidence, according to our model.

$$\begin{aligned} \text{conclusion} &= \arg \max_{\text{possible truth}} P(\text{possible truth}|\text{evidence}) \\ &\stackrel{\square}{=} \arg \max_{\text{possible truth}} \frac{P(\text{evidence}|\text{possible truth})P(\text{possible truth})}{P(\text{evidence})} \end{aligned}$$

In the above equation, we directly apply the Bayes’ theorem.

$$\begin{aligned} \text{conclusion} &= \arg \max_{\text{possible truth}} P(\text{possible truth}|\text{evidence}) \\ &= \arg \max_{\text{possible truth}} \frac{P(\text{evidence}|\text{possible truth})P(\text{possible truth})}{P(\text{evidence})} \\ &\stackrel{\square}{=} \arg \max_{\text{possible truth}} P(\text{evidence}|\text{possible truth})P(\text{possible truth}) \end{aligned}$$

It is important to note that the reason why the above equation holds is because the denominator $P(\text{evidence})$ does not depend on ‘possible truth’; i.e., it is constant for different values of ‘possible truth’. A frequent mistake is to assume that the equation is true because the $P(\text{possible truth})$ is 1.

If you are not clear why ‘evidence’ does not depend on different ‘possible truths’, it will be more clear once we express this in terms of random variables.

Bayesian Inference: Speech Recognition Example

Let us look at the speech recognition as an example of the problem that can be addressed using Bayesian Inference on a generative model.

Slide notes:

Bayesian Inference: Speech Recognition Example

- evidence \rightarrow sound
- possible truth \rightarrow utterance (words spoken)
- our best guess about utterance \rightarrow utterance*

$$\begin{aligned} \text{utterance}^* &= \arg \max_{\text{all utterances}} P(\text{utterance}|\text{sound}) \\ &= \arg \max_{\text{all utterances}} \frac{P(\text{sound}|\text{utterance})P(\text{utterance})}{P(\text{sound})} \\ &= \arg \max_{\text{utterance}} P(\text{sound}|\text{utterance})P(\text{utterance}) \end{aligned}$$

At the end, we need to estimate two probabilities to recognize the utterance base on the sound: $P(\text{sound}|\text{utterance})$ which is estimated using an *acoustic model*, and $P(\text{utterance})$ which is estimated using a *language model*.

11.4 Probabilistic Modeling

Slide notes:

Probabilistic Modeling

- How do we create and use a probabilistic model?
- Model elements:
 - Random variables
 - Model configuration (Random configuration)
 - Variable dependencies
 - Model parameters
- Computational tasks

Random Variables

Slide notes:

Random Variables

- Random variable V , defining an event as $V = x$ for some value x from a domain of values D ; i.e., $x \in D$
- $V = x$ is usually not a **basic** event due to having more variables
- An event with two random variables: $V_1 = x_1, V_2 = x_2$
- Multiple random variables: $\mathbf{V} = (V_1, V_2, \dots, V_n)$

It is frequently convenient to define events as $V = x$, where V is a variable and x is a value from a domain of values for that variable. This event can be a basic event, but usually it is not since we may have more variables.

Slide notes:

Model Configuration (Random Configuration)

- **Full Configuration:** If a model has n random variables, then a Full Model Configuration is an assignment of all the variables:

$$V_1 = x_1, V_2 = x_2, \dots, V_n = x_n$$

- **Partial configuration:** only some variables are assigned, e.g.:

$$V_1 = x_1, V_2 = x_2, \dots, V_k = x_k \quad (k < n)$$

Random Configuration. A **probabilistic model** includes a set of random variables with their domains, i.e., sets of the possible values of those variables. If all variables are assigned some values, we will call it a **random configuration** of the model, or a **full random configuration**. For example, if V_1, V_2, \dots, V_n are all random variables of a model, then an assignment of those variables, such as:

$$V_1 = x_1, V_2 = x_2, \dots, V_n = x_n$$

is a **full random configuration**. For each such configuration, the model provides a way of calculating its probability; i.e., calculating the value

$$P(V_1 = x_1, V_2 = x_2, \dots, V_n = x_n).$$

From the probability theory perspective, a model defines an event space, where each full random configuration is a simple event.

If we assume that random variables are ordered in a certain way, we do not have to list variable names but just the values, so each full random configuration becomes an n -tuple of values (x_1, x_2, \dots, x_n) , which we also call a **vector**, and use notation $\mathbf{x} = (x_1, x_2, \dots, x_n)$. Vectors are normally n -tuples of numbers, but we generalize this concept here. When creating a probabilistic model for a problem, we assume that a sequence of full configurations $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ is drawn from some random source:

$$\begin{aligned} \mathbf{x}^{(1)} &= (x_{11}, x_{12}, \dots, x_{1n}) \\ \mathbf{x}^{(2)} &= (x_{21}, x_{22}, \dots, x_{2n}) \\ &\vdots \\ \mathbf{x}^{(t)} &= (x_{t1}, x_{t2}, \dots, x_{tn}) \\ &\vdots \end{aligned}$$

Again, we assume a fixed number n of components in each configuration, and assume values x_{ij} are from a finite set of values that can be assigned to the variable V_j . For simplicity reasons, we will frequently assume that all variables get values from the same set $\{x_1, x_2, \dots, x_m\}$.

If only some variables are assigned, we will call such assignment a **partial configuration**. For example, if only variables V_1, V_2, \dots, V_k , where $k < n$ are assigned, then we have a partial configuration: $V_1 = x_1, V_2 = x_2, \dots, V_k = x_k$. In probabilistic terms, a partial configuration defines a random event that consists of all full configurations that satisfy the partial configuration assignment. This event can be described as the following set of full configurations:

$$\{(V_1 = y_1, V_2 = y_2, \dots, V_n = y_n) \mid y_1 = x_1 \wedge y_2 = x_2 \wedge \dots \wedge y_k = x_k\}$$

or, it is a set of all full configurations $(V_1 = y_1, V_2 = y_2, \dots, V_n = y_n)$ such that $y_1 = x_1$ and $y_2 = x_2$ and \dots and $y_k = x_k$.

Probabilistic modelling in NLP can be described as a general framework for modeling NLP problems using random variables, random configurations, and finding effective ways of reasoning about probabilities of these configurations.

Variable Independence and Dependence

- Random variables V_1 and V_2 are *independent* if $P(V_1 = x_1, V_2 = x_2) = P(V_1 = x_1)P(V_2 = x_2)$ for all x_1, x_2
- or expressed in a different way: $P(V_1 = x_1 | V_2 = x_2) = P(V_1 = x_1)$ for all x_1, x_2, x_3 .
- Random variables V_1 and V_2 are *conditionally independent given V_3* if, for all x_1, x_2, x_3 :
 - $P(V_1 = x_1, V_2 = x_2 | V_3 = x_3) =$
 - $P(V_1 = x_1 | V_3 = x_3)P(V_2 = x_2 | V_3 = x_3)$
- or
 - $P(V_1 = x_1 | V_2 = x_2, V_3 = x_3) = P(V_1 = x_1 | V_3 = x_3)$

As we will see, effective calculation of probabilities and probabilistic inference in general is based on making assumptions about dependence and independence of variables. Two random variables V_1 and V_2 are independent if:

$$P(V_1 = x_1, V_2 = x_2) = P(V_1 = x_1) \cdot P(V_2 = x_2)$$

11.5 Computational Tasks in Probabilistic Modeling

The computational tasks in probabilistic modeling can be divided into the following types of tasks:

- 1. Evaluation:** compute probability of a complete configuration
- 2. Simulation:** (sampling, generation) generate random configurations
- 3. Inference:** has the following sub-tasks:
 - 3.a Marginalization:** computing probability of a partial configuration,
 - 3.b Conditioning:** computing conditional probability of a completion given an observation,
 - 3.c Completion:** finding the most probable completion, given an observation
- 4. Learning:** learning parameters of a model from data.

We have four main tasks, and the task of probabilistic inference is divided further in to three subtasks. A more detailed description of these tasks follows:

1. Evaluation is the task of computing the probability of a full configuration in a probabilistic model; i.e., calculating

$$P(V_1 = x_1, V_2 = x_2, \dots, V_n = x_n)$$

A model definition usually provides a straightforward approach to calculate this probability given that we know all parameters of the model.

2. Simulation is the task of producing full configurations according to a given model. Those configurations should satisfy the model conditions, which typically means that probabilities of some events should converge to probabilities specified by the model over many simulated configurations.

3. Inference is divided into three different subtasks: marginalization, conditioning, and completion. **Marginalization (3.a)** is the problem of computing a marginal probability; i.e., the probability of a partial configuration, such as:

$$P(V_1 = x_1, V_2 = x_2, \dots, V_k = x_k),$$

where $1 \leq k < n$. **Conditioning (3.b)** is the problem of computing a conditional probability, which in terms of variables of a model typically means a conditional probability of the values of some variables, given the values of some other variables. For example, a conditioning task could be expressed as:

$$P(V_1 = x_1, V_2 = x_2, \dots, V_k = x_k | V_{k+1} = x_{k+1}, \dots, V_n = x_n).$$

Generally, we can do efficiently marginalization, we can also do conditioning by definition, as shown in this formula:

$$P(V_1 = x_1, V_2 = x_2, \dots, V_k = x_k \mid V_{k+1} = x_{k+1}, \dots, V_n = x_n) = \frac{P(V_1 = x_1, V_2 = x_2, \dots, V_k = x_k)}{P(V_{k+1} = x_{k+1}, \dots, V_n = x_n)}$$

Completion (3.c) is the problem of finding the most probable assignment of some variables, usually called **hidden variables**, given some other variables, usually called **observed variables**. We approach it by finding the assignments of the hidden variables for which the conditional probability of those assignments given assignments of observed variables is maximized. For example, this task can be described using the formula

$$\arg \max_{x_1, \dots, x_k} P(V_1 = x_1, V_2 = x_2, \dots, V_k = x_k \mid V_{k+1} = x_{k+1}, \dots, V_n = x_n).$$

In the completion formula $P(\alpha|\beta)$, the variables in the α part are called **hidden variables** and the variables in the β part are called **observed variables**, or **observables**, because normally know, or can observe, variables in the part β , and we do not know the hidden variables in the part α .

4. Learning is the task of determining the parameters of a probabilistic model given the data that the model is supposed to describe. Learning model parameters from incomplete data can be challenging. If we have a dataset with a list of full configurations we can calculate the parameters by counting the occurrences of interesting events. This is called Maximum Likelihood Estimation (MLE), since it can be shown that in this way we obtain a model that gives highest likelihood of data used in learning.

To illustrate probabilistic modelling and the computational tasks, we will use the following illustrative example in spam detection.

Illustrative Example: Spam Detection

The problem of spam detection in e-mail is the problem of automatically detecting whether an arbitrary e-mail message is spam or not. In a toy model, we assume that we can detect whether a message is spam or not relying only on the fact whether the ‘Subject:’ header of the message is capitalized (i.e., completely written in uppercase letters) and whether the ‘Subject:’ header contains the word ‘free’ (uppercase or lowercase). For example, “NEW MORTGAGE RATE” is likely the subject of a spam message, as well as “Money for Free,” “FREE lunch,” etc.

Hence, our model is based on the following three random variables and each of them gets one of two values Y (for Yes) or N (for No):

- Caps* = ‘Y’ if the message subject line does not contain lowercase letter, ‘N’ otherwise,
- Free* = ‘Y’ if the word ‘free’ appears in the message subject line (letter case is ignored),
‘N’ otherwise, and
- Spam* = ‘Y’ if the message is spam, and ‘N’ otherwise.

In order to learn what happens in real-world, we open our mailbox, which serves as our random source, randomly select 100 messages and count how many times each configuration appears.

We might obtain the following table:

<i>Free</i>	<i>Caps</i>	<i>Spam</i>	Number of messages
Y	Y	Y	20
Y	Y	N	1
Y	N	Y	5
Y	N	N	0
N	Y	Y	20
N	Y	N	3
N	N	Y	2
N	N	N	49
Total:			100

What are examples of computational tasks in this example?

Let us consider our first, straightforward model, called **Joint Distribution Model**.